Model Assumptions for MLR

Testing Model Fit 00000

Inference for Multiple Linear Regression

Nate Wells

Math 141, 4/20/22

	Regression

Model Assumptions for MLR 0000 Testing Model Fit

Model Selection

Outline

In this lecture, we will...

- Review framework for multilinear regression
- Discuss inference procedures for MLR models
- Investigate tools for "Model Selection"

Section 1

Multiple Linear Regression

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	
00000	0000	00000	00000000

 In a multiple linear regression model (MLR), we express the response variable Y as a linear combination of k explanatory variables X₁, X₂,..., X_k:

$$\hat{Y} = eta_0 + eta_1 \cdot X_1 + eta_2 \cdot X_2 + \dots + eta_k \cdot X_k$$

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	
00000	0000	00000	00000000

 In a multiple linear regression model (MLR), we express the response variable Y as a linear combination of k explanatory variables X₁, X₂,..., X_k:

$$\hat{Y} = \beta_0 + \beta_1 \cdot X_1 + \beta_2 \cdot X_2 + \cdots + \beta_k \cdot X_k$$

 We use the following R code to fit and summarize a linear model: mod<-lm(Y ~ X1 + X2 + X3, data = my_data) get_regression_table(mod)

	A CIDDIE.	4 X /					
	term	estimate	std_error	statistic	p_value	lower_ci	upper_ci
	<chr></chr>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>
1	intercept	3.26	7.94	0.41	0.686	-13.3	19.8
2	X1	-1.24	0.313	-3.95	0.001	-1.89	-0.584
3	X2	2.68	1.94	1.38	0.182	-1.36	6.72
4	X3	3.20	0.397	8.06	0	2.37	4.02
	" 1 2 3 4	<pre>term</pre>	term estimate <chr> <dbl> 1 intercept 3.26 2 X1 -1.24 3 X2 2.68 4 X3 3.20</dbl></chr>	*** *** **** term estimate std_error <chr> <dbl> <dbl> 1 intercept 3.26 7.94 2 X1 -1.24 0.313 3 X2 2.68 1.94 4 X3 3.20 0.397</dbl></dbl></chr>	term estimate std_error statistic <chr> <chr> <dbl></dbl> <dbl></dbl> 1 intercept 3.26 7.94 0.41 2 X1 -1.24 0.313 -3.95 3 X2 2.68 1.94 1.38 4 X3 3.20 0.397 8.06</chr></chr>	term estimate std_error statistic p_value <chr> <chr> <dbl><dbl><dbl><dbl><dbl><dbl><dbl><db< td=""><td>term estimate std_error statistic p_value lower_ci <chr> <chr> <dbl><dbl><dbl><dbl><dbl><dbl><dbl><db< td=""></db<></dbl></dbl></dbl></dbl></dbl></dbl></dbl></chr></chr></td></db<></dbl></dbl></dbl></dbl></dbl></dbl></dbl></chr></chr>	term estimate std_error statistic p_value lower_ci <chr> <chr> <dbl><dbl><dbl><dbl><dbl><dbl><dbl><db< td=""></db<></dbl></dbl></dbl></dbl></dbl></dbl></dbl></chr></chr>

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	
00000	0000	00000	00000000

 In a multiple linear regression model (MLR), we express the response variable Y as a linear combination of k explanatory variables X₁, X₂,..., X_k:

$$\hat{Y} = \beta_0 + \beta_1 \cdot X_1 + \beta_2 \cdot X_2 + \cdots + \beta_k \cdot X_k$$

 We use the following R code to fit and summarize a linear model: mod<-lm(Y ~ X1 + X2 + X3, data = my_data) get_regression_table(mod)

#	A tibble:	4 x 7					
	term	estimate	std_error	statistic	p_value	lower_ci	upper_ci
	<chr></chr>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>
1	intercept	3.26	7.94	0.41	0.686	-13.3	19.8
2	X1	-1.24	0.313	-3.95	0.001	-1.89	-0.584
3	X2	2.68	1.94	1.38	0.182	-1.36	6.72
4	X3	3.20	0.397	8.06	0	2.37	4.02
	# 1 2 3 4	<pre># A tibble: term <chr> 1 intercept 2 X1 3 X2 4 X3</chr></pre>	<pre># A tibble: 4 x 7 term estimate <chr></chr></pre>	<pre># A tibble: 4 x 7 term</pre>	<pre># A tibble: 4 x 7 term estimate std_error statistic <chr> <dbl> <dbl> <dbl> <dbl></dbl></dbl></dbl></dbl></chr></pre>	<pre># A tibble: 4 x 7 term estimate std_error statistic p_value <chr> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl></dbl></dbl></dbl></dbl></dbl></dbl></chr></pre>	<pre># A tibble: 4 x 7 term estimate std_error statistic p_value lower_ci <chr> <dbl> <dbl>dbl > <dbl>dbl >dbl >dbl >dbl >dbl >dbl >dbl ></dbl></dbl></dbl></dbl></dbl></dbl></dbl></dbl></dbl></dbl></dbl></dbl></dbl></dbl></dbl></dbl></dbl></dbl></dbl></dbl></dbl></dbl></dbl></dbl></dbl></dbl></dbl></dbl></dbl></dbl></dbl></dbl></dbl></dbl></dbl></dbl></dbl></dbl></dbl></dbl></dbl></dbl></dbl></dbl></dbl></dbl></dbl></dbl></dbl></dbl></dbl></dbl></dbl></dbl></chr></pre>

Which gives us our linear regression formula:

 $\hat{Y} = 3.26 - 1.24 \cdot X_1 + 2.68 \cdot X_2 + 3.2 \cdot X_3$

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	
00000	0000	00000	00000000

 In a multiple linear regression model (MLR), we express the response variable Y as a linear combination of k explanatory variables X₁, X₂,..., X_k:

$$\hat{Y} = \beta_0 + \beta_1 \cdot X_1 + \beta_2 \cdot X_2 + \cdots + \beta_k \cdot X_k$$

• We use the following R code to fit and summarize a linear model: mod<-lm(Y ~ X1 + X2 + X3, data = my_data) get_regression_table(mod)

##	#	A tibble:	4 x 7					
##		term	estimate	std_error	statistic	p_value	lower_ci	upper_ci
##		<chr></chr>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>
##	1	intercept	3.26	7.94	0.41	0.686	-13.3	19.8
##	2	X1	-1.24	0.313	-3.95	0.001	-1.89	-0.584
##	3	X2	2.68	1.94	1.38	0.182	-1.36	6.72
##	4	ХЗ	3.20	0.397	8.06	0	2.37	4.02

Which gives us our linear regression formula:

 $\hat{Y} = 3.26 - 1.24 \cdot X_1 + 2.68 \cdot X_2 + 3.2 \cdot X_3$

• The slope on each variable indicates the changed in the predicted value of Y per unit change in that variable, with all other variables held constant

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	
00000			

• How does page count vary across year of publication and division?

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	
00000			

• How does page count vary across year of publication and division?



Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	
000000			

• How does page count vary across year of publication and division?

Multiple Linear Regression	Model Assumptions for MLR 0000	Testing Model Fit 00000	Model Selection

.

```
    How does page count vary across year of publication and division?
theses_mlr <- lm(n_pages ~ year + division, data = theses)
get_regression_table(theses_mlr)
```

##	#	A tibble: /	X /					
##		term	estimate	std_error	statistic	p_value	lower_ci	upper_ci
##		<chr></chr>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>
##	1	intercept	1149.	126.	9.14	0	902.	1396.
##	2	year	-0.54	0.063	-8.58	0	-0.664	-0.416
##	3	divisionHSS	34.5	5.61	6.16	0	23.5	45.5
##	4	divisionID	18.7	7.38	2.53	0.011	4.22	33.2
##	5	divisionLL	12.6	5.84	2.16	0.031	1.16	24.1
##	6	divisionMNS	-11.2	5.50	-2.03	0.042	-22.0	-0.389
##	7	divisionPRPI	8.28	5.92	1.40	0.162	-3.34	19.9

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	Model Selection
000●00	0000	00000	
Reed Theses			

How does page count vary across year of publication and division?
 theses_mlr <- lm(n_pages ~ year + division, data = theses)
 get_regression_table(theses_mlr)

A tibble: 7 x 7 ## term estimate std error statistic p value lower ci upper ci <chr>> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> ## 126. ## 1 intercept 1149. 9.14 0 902. 1396. ## 2 year -0.664 -0.54 0.063 -8.58 0 -0.416## 3 divisionHSS 34.5 5.61 6.16 23.5 45.5 0 18.7 7.38 2.53 0.011 4.22 33.2 ## 4 divisionID ## 5 divisionLL 12.6 5.84 2.16 0.031 1.16 24.1 ## 6 divisionMNS -11.2 5.50 -2.03 0.042 -22.0 -0.389 ## 7 divisionPRPL 8.28 5.92 1.40 0.162 -3.34 19.9

 $\mathrm{Pages} = 1149 - 0.54 \cdot \mathrm{Year} + 34.5 \cdot \mathrm{HSS} + 18.7 \cdot \mathrm{ID} + 12.6 \cdot \mathrm{LL} - 11.2 \cdot \mathrm{MNS} + 8.3 \cdot \mathrm{PRPL}$

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	Model Selection
000●00	0000	00000	
Reed Theses			

 How does page count vary across year of publication and division? theses_mlr <- lm(n_pages ~ year + division, data = theses) get_regression_table(theses_mlr)

A tibble: 7 x 7

##		term	estimate	std_error	statistic	p_value	lower_ci	upper_ci
##		<chr></chr>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>
##	1	intercept	1149.	126.	9.14	0	902.	1396.
##	2	year	-0.54	0.063	-8.58	0	-0.664	-0.416
##	3	divisionHSS	34.5	5.61	6.16	0	23.5	45.5
##	4	divisionID	18.7	7.38	2.53	0.011	4.22	33.2
##	5	divisionLL	12.6	5.84	2.16	0.031	1.16	24.1
##	6	divisionMNS	-11.2	5.50	-2.03	0.042	-22.0	-0.389
##	7	divisionPRPL	8.28	5.92	1.40	0.162	-3.34	19.9

 $\mathrm{Pages} = 1149 - 0.54 \cdot \mathrm{Year} + 34.5 \cdot \mathrm{HSS} + 18.7 \cdot \mathrm{ID} + 12.6 \cdot \mathrm{LL} - 11.2 \cdot \mathrm{MNS} + 8.3 \cdot \mathrm{PRPL}$

• Which division is used as the baseline? What does the intercept represent?

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	Model Selection
000●00	0000	00000	

```
    How does page count vary across year of publication and division?
    theses_mlr <- lm(n_pages ~ year + division, data = theses)</li>
    get_regression_table(theses_mlr)
```

A tibble: 7 x 7

##		term	estimate	std_error	statistic	p_value	lower_ci	upper_ci
##		<chr></chr>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>
##	1	intercept	1149.	126.	9.14	0	902.	1396.
##	2	year	-0.54	0.063	-8.58	0	-0.664	-0.416
##	3	divisionHSS	34.5	5.61	6.16	0	23.5	45.5
##	4	divisionID	18.7	7.38	2.53	0.011	4.22	33.2
##	5	divisionLL	12.6	5.84	2.16	0.031	1.16	24.1
##	6	divisionMNS	-11.2	5.50	-2.03	0.042	-22.0	-0.389
##	7	divisionPRPL	8.28	5.92	1.40	0.162	-3.34	19.9

 $\mathrm{Pages} = 1149 - 0.54 \cdot \mathrm{Year} + 34.5 \cdot \mathrm{HSS} + 18.7 \cdot \mathrm{ID} + 12.6 \cdot \mathrm{LL} - 11.2 \cdot \mathrm{MNS} + 8.3 \cdot \mathrm{PRPL}$

- Which division is used as the baseline? What does the intercept represent?
- What does the coefficient of 34.5 on HSS mean in context?

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	Model Selection
000●00	0000	00000	

```
    How does page count vary across year of publication and division?
    theses_mlr <- lm(n_pages ~ year + division, data = theses)</li>
    get_regression_table(theses_mlr)
```

A tibble: 7 x 7

##		term	estimate	std_error	statistic	p_value	lower_ci	upper_ci
##		<chr></chr>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>
##	1	intercept	1149.	126.	9.14	0	902.	1396.
##	2	year	-0.54	0.063	-8.58	0	-0.664	-0.416
##	3	divisionHSS	34.5	5.61	6.16	0	23.5	45.5
##	4	divisionID	18.7	7.38	2.53	0.011	4.22	33.2
##	5	divisionLL	12.6	5.84	2.16	0.031	1.16	24.1
##	6	divisionMNS	-11.2	5.50	-2.03	0.042	-22.0	-0.389
##	7	divisionPRPL	8.28	5.92	1.40	0.162	-3.34	19.9

 $\mathrm{Pages} = 1149 - 0.54 \cdot \mathrm{Year} + 34.5 \cdot \mathrm{HSS} + 18.7 \cdot \mathrm{ID} + 12.6 \cdot \mathrm{LL} - 11.2 \cdot \mathrm{MNS} + 8.3 \cdot \mathrm{PRPL}$

- Which division is used as the baseline? What does the intercept represent?
- What does the coefficient of 34.5 on HSS mean in context?
- What does the coefficient of -0.54 on Year mean in context?

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	
000000			

• The regression table provides p-values for each slope in the model.

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	
000000			

- The regression table provides p-values for each slope in the model.
 - But what hypotheses are being tested?

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	
000000			

- The regression table provides p-values for each slope in the model.
 - But what hypotheses are being tested?
- In a **MLR model**, we are still interested in determining whether a slope β_i is 0.

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	
000000			

- The regression table provides p-values for each slope in the model.
 - But what hypotheses are being tested?
- In a **MLR model**, we are still interested in determining whether a slope β_i is 0.
 - But we want to investigate this slope in light of the other variables in the model.

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	
000000			

- The regression table provides p-values for each slope in the model.
 - But what hypotheses are being tested?
- In a **MLR model**, we are still interested in determining whether a slope β_i is 0.
 - But we want to investigate this slope in light of the other variables in the model.
- Each row corresponds to a hypothesis test of the form

 $H_0: \beta_i = 0$, given that other variables are included in the model

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	
000000			

- The regression table provides p-values for each slope in the model.
 - But what hypotheses are being tested?
- In a **MLR model**, we are still interested in determining whether a slope β_i is 0.
 - But we want to investigate this slope in light of the other variables in the model.
- Each row corresponds to a hypothesis test of the form

 $H_0: \beta_i = 0$, given that other variables are included in the model

• I.e. The year row corresponds to the test of

 $H_0: \beta_{year} = 0$, given each division is included in the model

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	
000000			

- The regression table provides p-values for each slope in the model.
 - But what hypotheses are being tested?
- In a **MLR model**, we are still interested in determining whether a slope β_i is 0.
 - But we want to investigate this slope in light of the other variables in the model.
- Each row corresponds to a hypothesis test of the form

 $H_0: \beta_i = 0$, given that other variables are included in the model

• I.e. The year row corresponds to the test of

 $H_0: \beta_{year} = 0$, given each division is included in the model

• Reminder: The p-value is the probability of obtaining a statistic as extreme as the observed statistic, if the null hypothesis were true.

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	
000000			

- The regression table provides p-values for each slope in the model.
 - But what hypotheses are being tested?
- In a **MLR model**, we are still interested in determining whether a slope β_i is 0.
 - But we want to investigate this slope in light of the other variables in the model.
- Each row corresponds to a hypothesis test of the form

 $H_0: \beta_i = 0$, given that other variables are included in the model

• I.e. The year row corresponds to the test of

 $H_0: \beta_{year} = 0$, given each division is included in the model

- Reminder: The p-value is the probability of obtaining a statistic as extreme as the observed statistic, **if the null hypothesis were true**.
- The standard error, statistic, and p-values are all calculated using theory-based methods.

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	
000000			

- The regression table provides p-values for each slope in the model.
 - But what hypotheses are being tested?
- In a **MLR model**, we are still interested in determining whether a slope β_i is 0.
 - But we want to investigate this slope in light of the other variables in the model.
- Each row corresponds to a hypothesis test of the form

 $H_0: \beta_i = 0$, given that other variables are included in the model

• I.e. The year row corresponds to the test of

 $H_0: \beta_{year} = 0$, given each division is included in the model

- Reminder: The p-value is the probability of obtaining a statistic as extreme as the observed statistic, **if the null hypothesis were true**.
- The standard error, statistic, and p-values are all calculated using theory-based methods.
 - But the formula is very complicated, requiring matrices and linear algebra (If interested, take Math 392)

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	
00000			

##	Ŧ	A tibble: /	X /					
##		term	estimate	std_error	statistic	p_value	lower_ci	upper_ci
##		<chr></chr>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>
##	1	intercept	1149.	126.	9.14	0	902.	1396.
##	2	year	-0.54	0.063	-8.58	0	-0.664	-0.416
##	3	divisionHSS	34.5	5.61	6.16	0	23.5	45.5
##	4	divisionID	18.7	7.38	2.53	0.011	4.22	33.2
##	5	divisionLL	12.6	5.84	2.16	0.031	1.16	24.1
##	6	divisionMNS	-11.2	5.50	-2.03	0.042	-22.0	-0.389
##	7	divisionPRPI	8.28	5.92	1.40	0.162	-3.34	19.9

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	
00000			

• Consider the regression table... -

##	#	A tibble: 7	x 7					
##		term	estimate	std_error	statistic	p_value	lower_ci	upper_ci
##		<chr></chr>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>
##	1	intercept	1149.	126.	9.14	0	902.	1396.
##	2	year	-0.54	0.063	-8.58	0	-0.664	-0.416
##	3	divisionHSS	34.5	5.61	6.16	0	23.5	45.5
##	4	divisionID	18.7	7.38	2.53	0.011	4.22	33.2
##	5	divisionLL	12.6	5.84	2.16	0.031	1.16	24.1
##	6	divisionMNS	-11.2	5.50	-2.03	0.042	-22.0	-0.389
##	7	divisionPRPI	. 8.28	5.92	1.40	0.162	-3.34	19.9

For which hypothesis tests would you reject H_0 ? ٠

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	
00000			

....

##	Ŧ	A tibble: /	X /					
##		term	estimate	std_error	statistic	p_value	lower_ci	upper_ci
##		<chr></chr>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>
##	1	intercept	1149.	126.	9.14	0	902.	1396.
##	2	year	-0.54	0.063	-8.58	0	-0.664	-0.416
##	3	divisionHSS	34.5	5.61	6.16	0	23.5	45.5
##	4	divisionID	18.7	7.38	2.53	0.011	4.22	33.2
##	5	divisionLL	12.6	5.84	2.16	0.031	1.16	24.1
##	6	divisionMNS	-11.2	5.50	-2.03	0.042	-22.0	-0.389
##	7	divisionPRPI	. 8.28	5.92	1.40	0.162	-3.34	19.9

- For which hypothesis tests would you reject *H*₀?
 - What does this mean in context?

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	
00000			

##	#	A tibble: 7	x 7					
##		term	estimate	std_error	statistic	p_value	lower_ci	upper_ci
##		<chr></chr>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>
##	1	intercept	1149.	126.	9.14	0	902.	1396.
##	2	year	-0.54	0.063	-8.58	0	-0.664	-0.416
##	3	divisionHSS	34.5	5.61	6.16	0	23.5	45.5
##	4	divisionID	18.7	7.38	2.53	0.011	4.22	33.2
##	5	divisionLL	12.6	5.84	2.16	0.031	1.16	24.1
##	6	divisionMNS	-11.2	5.50	-2.03	0.042	-22.0	-0.389
##	7	divisionPRPI	. 8.28	5.92	1.40	0.162	-3.34	19.9

- For which hypothesis tests would you reject *H*₀?
 - What does this mean in context?
- For which would you fail to reject H_0 ?

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	
00000			

##	#	A tibble: 7	x 7					
##		term	estimate	std_error	statistic	p_value	lower_ci	upper_ci
##		<chr></chr>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>
##	1	intercept	1149.	126.	9.14	0	902.	1396.
##	2	year	-0.54	0.063	-8.58	0	-0.664	-0.416
##	3	divisionHSS	34.5	5.61	6.16	0	23.5	45.5
##	4	divisionID	18.7	7.38	2.53	0.011	4.22	33.2
##	5	divisionLL	12.6	5.84	2.16	0.031	1.16	24.1
##	6	divisionMNS	-11.2	5.50	-2.03	0.042	-22.0	-0.389
##	7	divisionPRPI	. 8.28	5.92	1.40	0.162	-3.34	19.9

- For which hypothesis tests would you reject *H*₀?
 - What does this mean in context?
- For which would you fail to reject H₀?
 - What does this mean in context?

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	
00000			

##	Ŧ	A tibble: /	X /					
##		term	estimate	std_error	statistic	p_value	lower_ci	upper_ci
##		<chr></chr>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>
##	1	intercept	1149.	126.	9.14	0	902.	1396.
##	2	year	-0.54	0.063	-8.58	0	-0.664	-0.416
##	3	divisionHSS	34.5	5.61	6.16	0	23.5	45.5
##	4	divisionID	18.7	7.38	2.53	0.011	4.22	33.2
##	5	divisionLL	12.6	5.84	2.16	0.031	1.16	24.1
##	6	divisionMNS	-11.2	5.50	-2.03	0.042	-22.0	-0.389
##	7	divisionPRPI	. 8.28	5.92	1.40	0.162	-3.34	19.9

- For which hypothesis tests would you reject *H*₀?
 - What does this mean in context?
- For which would you fail to reject H₀?
 - What does this mean in context?
- What is the relationship between the estimate, the std_error and lower_ci, upper_ci?

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	
00000			

##	#	A tibble: 7	x 7					
##		term	estimate	std_error	statistic	p_value	lower_ci	upper_ci
##		<chr></chr>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>
##	1	intercept	1149.	126.	9.14	0	902.	1396.
##	2	year	-0.54	0.063	-8.58	0	-0.664	-0.416
##	3	divisionHSS	34.5	5.61	6.16	0	23.5	45.5
##	4	divisionID	18.7	7.38	2.53	0.011	4.22	33.2
##	5	divisionLL	12.6	5.84	2.16	0.031	1.16	24.1
##	6	divisionMNS	-11.2	5.50	-2.03	0.042	-22.0	-0.389
##	7	divisionPRPI	. 8.28	5.92	1.40	0.162	-3.34	19.9

- For which hypothesis tests would you reject H_0 ?
 - What does this mean in context?
- For which would you fail to reject H₀?
 - What does this mean in context?
- What is the relationship between the estimate, the std_error and lower_ci, upper_ci?
- How does the coefficient on year in the MLR model compare to the coefficient on year in the SLR model?

Model Assumptions for MLR

Testing Model Fit 00000 Model Selection

Section 2

Model Assumptions for MLR

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	
	0000		

• In order to responsibly use MLR to make inference, we need...

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	
	0000		

- In order to responsibly use MLR to make inference, we need...
- The relationship between explanatory and response variables must be approximately multilinear linear. (Linear)
- **2** The observations should be independent of one another. (Independence)
- The distribution of residuals should be bell-shaped, unimodal, symmetric, and centered at 0. (Normal)
- The variability of residuals should be roughly constant across entire data set. (Equal Variability)

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	
	0000		

- In order to responsibly use MLR to make inference, we need...
- The relationship between explanatory and response variables must be approximately multilinear linear. (Linear)
- **2** The observations should be independent of one another. (Independence)
- The distribution of residuals should be bell-shaped, unimodal, symmetric, and centered at 0. (Normal)
- The variability of residuals should be roughly constant across entire data set. (Equal Variability)
- How do we check some of these conditions? Why can't we create a scatterplot of residuals as we did for SLR?

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	
	0000		

- In order to responsibly use MLR to make inference, we need...
- The relationship between explanatory and response variables must be approximately multilinear linear. (Linear)
- **2** The observations should be independent of one another. (Independence)
- The distribution of residuals should be bell-shaped, unimodal, symmetric, and centered at 0. (Normal)
- The variability of residuals should be roughly constant across entire data set. (Equal Variability)
- How do we check some of these conditions? Why can't we create a scatterplot of residuals as we did for SLR?
 - Instead, we will use a scatterplot of residuals vs predicted values
| Model Assumptions for MLR | Testing Model Fit | |
|---------------------------|-------------------|--|
| 0000 | | |

```
mod res <- get regression points(theses mlr)</pre>
mod res %>%
 select(n_pages, n_pages_hat, residual)
## # A tibble: 772 x 3
##
     n_pages n_pages_hat residual
##
       <dbl>
                   <dbl>
                           <dbl>
                   72.5 11.5
##
   1
          84
##
   2
         139
                    88.2 50.8
                    82.2 -32.2
##
   3
          50
                    53.6
                          25.4
##
   4
          79
##
   5
          36
                   47.1
                          -11.1
##
   6
          80
                    54.1
                          25.9
##
   7
         134
                   67.6
                          66.4
##
   8
          58
                   75.2 -17.2
##
          69
                   75.7
                          -6.74
   9
## 10
          74
                    85.4
                           -11.4
     ... with 762 more rows
##
  #
```

	Model Assumptions for MLR	Testing Model Fit	
000000	0000	00000	00000000

```
mod_res %>% ggplot(aes(x = n_pages_hat, y = residual)
mod res <- get regression points(theses mlr</pre>
                                                geom point()+
mod res %>%
                                                geom smooth(method = "lm", se = F)+
  select(n_pages, n_pages_hat, residual)
                                                labs(title = "Residual vs Predicted")
                                                        Residual vs Predicted
   # A tibble: 772 x 3
##
      n_pages n_pages_hat residual
##
                                                                 ...
                                                                                 :
##
        <dbl>
                     <dbl>
                             <dbl>
                                                                                       ٠.
                     72.5 11.5
##
    1
           84
                                                      10
##
    2
          139
                     88.2
                            50.8
                     82.2
                            -32.2
##
    3
           50
                     53.6
                            25.4
##
    4
           79
                                                    esidual
##
    5
           36
                     47.1
                            -11.1
##
    6
           80
                     54.1
                            25.9
##
    7
          134
                     67.6
                            66.4
##
    8
           58
                     75.2 -17.2
           69
                     75.7
                            -6.74
                                                      -50
##
    9
           74
                      85.4
                             -11.4
##
   10
         with 762 more rows
                                                                       n pages hat
```

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	
	0000		



When analyzing residual vs. predicted plots, look for...

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	
	0000		



- When analyzing residual vs. predicted plots, look for...
 - Non-linear patterns
 - Increasing variability across range of predicted values
 - Outliers with atypical predicted value or large residual

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	
	0000		

Distribution of Residuals

• We can still look at the histogram and QQ Plot of residuals, as we did for SLR:

	Model Assumptions for MLR	Testing Model Fit	
000000	0000	00000	00000000

Distribution of Residuals

• We can still look at the histogram and QQ Plot of residuals, as we did for SLR:



000000 000 00000 0000000	Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	
		0000		

Distribution of Residuals

• We can still look at the histogram and QQ Plot of residuals, as we did for SLR:



• We see some evidence that residuals are not Normally distributed (right-skew)

	Model Assumptions for MLR	Testing Model Fit	
000000	0000	● 0000	00000000

Section 3

Testing Model Fit

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	Model Selection
000000	0000	○●○○○	

• Previously, we performed individual hypothesis tests to measure each explanatory variable's contribution to the model.

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	
		00000	

- Previously, we performed individual hypothesis tests to measure each explanatory variable's contribution to the model.
 - But, when considering a model with many explanatory variables, we need to be cautious about interpreting individual p-values.

	Model Assumptions for MLR	Testing Model Fit	
000000	0000	0000	00000000

- Previously, we performed individual hypothesis tests to measure each explanatory variable's contribution to the model.
 - But, when considering a model with many explanatory variables, we need to be cautious about interpreting individual p-values.
 - Suppose we had a model which contained 100 explanatory variables that were all independent of the response.
 - Approximately how many of the rows would have a p-value that is significant at the 0.05 level?

	Model Assumptions for MLR	Testing Model Fit	
000000	0000	0000	00000000

- Previously, we performed individual hypothesis tests to measure each explanatory variable's contribution to the model.
 - But, when considering a model with many explanatory variables, we need to be cautious about interpreting individual p-values.
 - Suppose we had a model which contained 100 explanatory variables that were all independent of the response.
 - Approximately how many of the rows would have a p-value that is significant at the 0.05 level?
- So we may also be interested in assessing how well the overall model fits the data.

	Model Assumptions for MLR	Testing Model Fit	
000000	0000	0000	00000000

- Previously, we performed individual hypothesis tests to measure each explanatory variable's contribution to the model.
 - But, when considering a model with many explanatory variables, we need to be cautious about interpreting individual p-values.
 - Suppose we had a model which contained 100 explanatory variables that were all independent of the response.
 - Approximately how many of the rows would have a p-value that is significant at the 0.05 level?
- So we may also be interested in assessing how well the overall model fits the data.
- We will test the following hypotheses:

	Model Assumptions for MLR	Testing Model Fit	
000000	0000	0000	00000000

- Previously, we performed individual hypothesis tests to measure each explanatory variable's contribution to the model.
 - But, when considering a model with many explanatory variables, we need to be cautious about interpreting individual p-values.
 - Suppose we had a model which contained 100 explanatory variables that were all independent of the response.
 - Approximately how many of the rows would have a p-value that is significant at the 0.05 level?
- So we may also be interested in assessing how well the overall model fits the data.
- We will test the following hypotheses:

$$H_0: \beta_1 = \beta_2 = \cdots = \beta_k = 0$$
 $H_a:$ At least one $\beta_i \neq 0$

	Model Assumptions for MLR	Testing Model Fit	
000000	0000	0000	00000000

- Previously, we performed individual hypothesis tests to measure each explanatory variable's contribution to the model.
 - But, when considering a model with many explanatory variables, we need to be cautious about interpreting individual p-values.
 - Suppose we had a model which contained 100 explanatory variables that were all independent of the response.
 - Approximately how many of the rows would have a p-value that is significant at the 0.05 level?
- So we may also be interested in assessing how well the overall model fits the data.
- We will test the following hypotheses:

$$H_0: \beta_1 = \beta_2 = \cdots = \beta_k = 0$$
 $H_a:$ At least one $\beta_i \neq 0$

i.e. We will assess how likely it was to obtain slope estimates β₁,..., β_k as large as those observed, if all the explanatory variables were uncorrelated with the response.

	Model Assumptions for MLR	Testing Model Fit	
000000	0000	0000	00000000

- Previously, we performed individual hypothesis tests to measure each explanatory variable's contribution to the model.
 - But, when considering a model with many explanatory variables, we need to be cautious about interpreting individual p-values.
 - Suppose we had a model which contained 100 explanatory variables that were all independent of the response.
 - Approximately how many of the rows would have a p-value that is significant at the 0.05 level?
- So we may also be interested in assessing how well the overall model fits the data.
- We will test the following hypotheses:

$$H_0: \beta_1 = \beta_2 = \cdots = \beta_k = 0$$
 $H_a:$ At least one $\beta_i \neq 0$

- i.e. We will assess how likely it was to obtain slope estimates β₁,..., β_k as large as those observed, if all the explanatory variables were uncorrelated with the response.
 - Note that $\beta_0 = 0$ is not in the null hypothesis, since we are only looking for evidence that at least one explanatory variable contributes to the response.

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	
		00000	

• Recall that the general form of a MLR model is

$$\mathbf{Y} = \beta_0 + \beta_1 \cdot \mathbf{X}_1 + \beta_2 \cdot \mathbf{X}_2 + \dots + \beta_k \cdot \mathbf{X}_k + \epsilon$$

Model Assumptions for MLR	Testing Model Fit	
	00000	

Recall that the general form of a MLR model is

$$\mathbf{Y} = \beta_0 + \beta_1 \cdot \mathbf{X}_1 + \beta_2 \cdot \mathbf{X}_2 + \dots + \beta_k \cdot \mathbf{X}_k + \epsilon$$

- We can split the total variability in the response variable into two pieces:
 - One representing the variability explained by the model
 - · Another representing the variability given by the errors

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	
		00000	

Recall that the general form of a MLR model is

$$Y = \beta_0 + \beta_1 \cdot X_1 + \beta_2 \cdot X_2 + \dots + \beta_k \cdot X_k + \epsilon$$

- We can split the total variability in the response variable into two pieces:
 - One representing the variability explained by the model
 - · Another representing the variability given by the errors
- We write...

SSTotal = SSModel + SSE

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	
		00000	

Recall that the general form of a MLR model is

$$Y = \beta_0 + \beta_1 \cdot X_1 + \beta_2 \cdot X_2 + \dots + \beta_k \cdot X_k + \epsilon$$

- We can split the total variability in the response variable into two pieces:
 - One representing the variability explained by the model
 - Another representing the variability given by the errors
- We write...

$$SSTotal = SSModel + SSE$$

where...

$$\begin{split} \text{SSTotal} &= \text{Total Sum of Squares} = \sum (y_i - \bar{y})^2 \\ \text{SSModel} &= \text{Sum of Squares Explained by Model} = \sum (\hat{y}_i - \bar{y})^2 \\ \text{SSE} &= \text{Sum of Squared Error} = \sum (y_i - \hat{y}_i)^2 \end{split}$$

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	
		00000	

Recall that the general form of a MLR model is

$$\mathbf{Y} = \beta_0 + \beta_1 \cdot \mathbf{X}_1 + \beta_2 \cdot \mathbf{X}_2 + \dots + \beta_k \cdot \mathbf{X}_k + \epsilon$$

- We can split the total variability in the response variable into two pieces:
 - One representing the variability explained by the model
 - · Another representing the variability given by the errors
- We write...

$$SSTotal = SSModel + SSE$$

where...

$$\begin{split} \text{SSTotal} &= \text{Total Sum of Squares} = \sum (y_i - \bar{y})^2 \\ \text{SSModel} &= \text{Sum of Squares Explained by Model} = \sum (\hat{y}_i - \bar{y})^2 \\ \text{SSE} &= \text{Sum of Squared Error} = \sum (y_i - \hat{y}_i)^2 \end{split}$$

• Note: SSE was the quantity we sought to minimize when fitting the least squares line.

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	Model Selection
000000	0000	000●0	

• A "good" model is one where almost all of the variability in response is due to the model, rather than the error.

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	
		00000	

- A "good" model is one where almost all of the variability in response is due to the model, rather than the error.
- Consider the following ratio:

$$F = rac{ ext{SSModel}}{ ext{SSE}} \cdot rac{n-k-1}{k}$$

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	
		00000	

- A "good" model is one where almost all of the variability in response is due to the model, rather than the error.
- Consider the following ratio:

$$F = \frac{\text{SSModel}}{\text{SSE}} \cdot \frac{n-k-1}{k}$$

- Good models should have values of F (much) larger than 1
- Models with variables uncorrelated with response should have values of F close to 1.

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	
		00000	

- A "good" model is one where almost all of the variability in response is due to the model, rather than the error.
- Consider the following ratio:

$$F = \frac{\text{SSModel}}{\text{SSE}} \cdot \frac{n-k-1}{k}$$

- Good models should have values of F (much) larger than 1
- Models with variables uncorrelated with response should have values of F close to 1.
- The extra $\frac{n-k-1}{k}$ eliminates bias due to the number of variables and observations

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	
		00000	

- A "good" model is one where almost all of the variability in response is due to the model, rather than the error.
- Consider the following ratio:

$$F = \frac{\text{SSModel}}{\text{SSE}} \cdot \frac{n-k-1}{k}$$

- Good models should have values of F (much) larger than 1
- Models with variables uncorrelated with response should have values of F close to 1.
- The extra $\frac{n-k-1}{k}$ eliminates bias due to the number of variables and observations
- Under the null hypothesis, F follows a certain theoretical distribution
 - This distribution is called the F distribution, which has two parameters: (k, n k 1)



Multiple Linear Regression 000000	Model Assumptions for MLR	Testing Model Fit 0000●	Model Selection

• We **could** compute the value of the F statistic by hand, and then compute the p-value of the F stat using an appropriate function in R (like pf).

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit 0000●	Model Selection

- We **could** compute the value of the F statistic by hand, and then compute the p-value of the F stat using an appropriate function in R (like pf).
 - BUT! The 1m function already does all of this for us.

Multiple Linear Regression 000000	Model Assumptions for MLR	Testing Model Fit 0000●	Model Selection

• We **could** compute the value of the F statistic by hand, and then compute the p-value of the F stat using an appropriate function in R (like pf).

BUT! The lm function already does all of this for us.

```
theses_mlr <- lm(n_pages ~ year + division, data = theses)
get_regression_summaries(theses_mlr)</pre>
```

A tibble: 1 x 9
r_squared adj_r_squared mse rmse sigma statistic p_value df nobs
<dbl> <dbl < dbl < dbl

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit 0000●	Model Selection

• We **could** compute the value of the F statistic by hand, and then compute the p-value of the F stat using an appropriate function in R (like pf).

• BUT! The 1m function already does all of this for us.

```
theses_mlr <- lm(n_pages ~ year + division, data = theses)
get_regression_summaries(theses_mlr)</pre>
```

A tibble: 1 x 9
r_squared adj_r_squared mse rmse sigma statistic p_value df nobs
<dbl> <dbl <dbl <dbl > <dbl

• The statistic in the get_regression_summaries table IS the F-statistic.

Multiple Linear Regression 000000	Model Assumptions for MLR	Testing Model Fit 0000●	Model Selection

• We **could** compute the value of the F statistic by hand, and then compute the p-value of the F stat using an appropriate function in R (like pf).

BUT! The lm function already does all of this for us.
 theses_mlr <- lm(n_pages ~ year + division, data = theses)
 get regression summaries(theses mlr)

- The statistic in the get_regression_summaries table IS the F-statistic.
- The p_value is the p-value for this statistic in the appropriate F distribution.

Multiple Linear Regression 000000	Model Assumptions for MLR	Testing Model Fit 0000●	Model Selection

• We **could** compute the value of the F statistic by hand, and then compute the p-value of the F stat using an appropriate function in R (like pf).

• BUT! The lm function already does all of this for us. theses_mlr <- lm(n_pages - year + division, data = theses)

```
get_regression_summaries(theses_mlr)
```

##	#	A tibble:	1 x 9							
##		r_squared	adj_r_squared	mse	rmse	sigma	statistic	p_value	df	nobs
##		<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>
##	1	0.254	0.248	1319.	36.3	36.5	43.5	0	6	772

- The statistic in the get_regression_summaries table IS the F-statistic.
- The p_value is the p-value for this statistic in the appropriate F distribution.
- In the case of the Reed Theses, we see that the P-value for the F test is (very close to) 0, and so we reject the null hypothesis.

Multiple Linear Regression 000000	Model Assumptions for MLR	Testing Model Fit 0000●	Model Selection

• We **could** compute the value of the F statistic by hand, and then compute the p-value of the F stat using an appropriate function in R (like pf).

• BUT! The lm function already does all of this for us. theses_mlr <- lm(n_pages ~ year + division, data = theses)

get_regression_summaries(theses_mlr)

##	#	A tibble:	1 x 9							
##		r_squared	adj_r_squared	mse	rmse	sigma	statistic	p_value	df	nobs
##		<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>
##	1	0.254	0.248	1319.	36.3	36.5	43.5	0	6	772

- The statistic in the get_regression_summaries table IS the F-statistic.
- The p_value is the p-value for this statistic in the appropriate F distribution.
- In the case of the Reed Theses, we see that the P-value for the F test is (very close to) 0, and so we reject the null hypothesis.
 - This sample gives evidence that at least one of the coefficients in the model is non-zero.

Model Assumptions for MLR

Testing Model Fit 00000 Model Selection

Section 4

Model Selection

Nate Wells

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	Model Selection
000000	0000	00000	○●○○○○○○

R^2 and Adjusted R^2

• In addition to the F statistic, we've already seen another metric from the goodness-of-fit of a multilinear model: R^2

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	Model Selection
			0000000

R^2 and Adjusted R^2

- In addition to the F statistic, we've already seen another metric from the goodness-of-fit of a multilinear model: R^2
- Recall:

$$R^2 = \frac{\text{SSModel}}{\text{TSS}} = \frac{\text{TSS} - \text{SSE}}{\text{TSS}} = 1 - \frac{\text{SSE}}{\text{TSS}}$$
Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	Model Selection
			0000000

- In addition to the F statistic, we've already seen another metric from the goodness-of-fit of a multilinear model: R^2
- Recall:

$$R^2 = \frac{\mathrm{SSModel}}{\mathrm{TSS}} = \frac{\mathrm{TSS} - \mathrm{SSE}}{\mathrm{TSS}} = 1 - \frac{\mathrm{SSE}}{\mathrm{TSS}}$$

• That is, R^2 measures the proportion of variability in the response explained by variability in explanatory variables.

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	Model Selection
000000	0000	00000	0000000

- In addition to the F statistic, we've already seen another metric from the goodness-of-fit of a multilinear model: R^2
- Recall:

$$R^2 = \frac{\text{SSModel}}{\text{TSS}} = \frac{\text{TSS} - \text{SSE}}{\text{TSS}} = 1 - \frac{\text{SSE}}{\text{TSS}}$$

- That is, R^2 measures the proportion of variability in the response explained by variability in explanatory variables.
 - If $R^2 \approx 1$, most of the variability in response is explained by linear relationship with explanatory variables.

	Model Assumptions for MLR	Testing Model Fit	Model Selection
000000	0000	00000	0000000

- In addition to the F statistic, we've already seen another metric from the goodness-of-fit of a multilinear model: R^2
- Recall:

$$R^2 = \frac{\text{SSModel}}{\text{TSS}} = \frac{\text{TSS} - \text{SSE}}{\text{TSS}} = 1 - \frac{\text{SSE}}{\text{TSS}}$$

- That is, R^2 measures the proportion of variability in the response explained by variability in explanatory variables.
 - If $R^2\approx 1,$ most of the variability in response is explained by linear relationship with explanatory variables.
 - While if $R^2 \approx 0$, almost none of the variability is explained by the model.

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	Model Selection
000000	0000	00000	0000000

- In addition to the F statistic, we've already seen another metric from the goodness-of-fit of a multilinear model: R^2
- Recall:

$$R^2 = \frac{\text{SSModel}}{\text{TSS}} = \frac{\text{TSS} - \text{SSE}}{\text{TSS}} = 1 - \frac{\text{SSE}}{\text{TSS}}$$

- That is, R^2 measures the proportion of variability in the response explained by variability in explanatory variables.
 - If $R^2\approx$ 1, most of the variability in response is explained by linear relationship with explanatory variables.
 - While if $R^2 \approx 0$, almost none of the variability is explained by the model.
- But it turns out this formula gives a **biased** estimate of the variability in the *population* explained by the model.

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	Model Selection
000000	0000	00000	0000000

- In addition to the F statistic, we've already seen another metric from the goodness-of-fit of a multilinear model: R^2
- Recall:

$$R^2 = \frac{\text{SSModel}}{\text{TSS}} = \frac{\text{TSS} - \text{SSE}}{\text{TSS}} = 1 - \frac{\text{SSE}}{\text{TSS}}$$

- That is, R^2 measures the proportion of variability in the response explained by variability in explanatory variables.
 - If $R^2\approx$ 1, most of the variability in response is explained by linear relationship with explanatory variables.
 - While if $R^2 \approx 0$, almost none of the variability is explained by the model.
- But it turns out this formula gives a **biased** estimate of the variability in the *population* explained by the model.
- Instead, we may use the adjusted R:

$$R_{\mathrm{adj}}^2 = 1 - \left(\frac{\mathrm{SSE}}{\mathrm{TSS}} \cdot \frac{n-1}{n-k-1}\right)$$

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	Model Selection
000000	0000	00000	0000000

- In addition to the F statistic, we've already seen another metric from the goodness-of-fit of a multilinear model: R^2
- Recall:

$$R^2 = rac{\mathrm{SSModel}}{\mathrm{TSS}} = rac{\mathrm{TSS} - \mathrm{SSE}}{\mathrm{TSS}} = 1 - rac{\mathrm{SSE}}{\mathrm{TSS}}$$

- That is, R^2 measures the proportion of variability in the response explained by variability in explanatory variables.
 - If $R^2\approx$ 1, most of the variability in response is explained by linear relationship with explanatory variables.
 - While if $R^2 \approx 0$, almost none of the variability is explained by the model.
- But it turns out this formula gives a **biased** estimate of the variability in the *population* explained by the model.
- Instead, we may use the adjusted R:

$$R_{\mathrm{adj}}^2 = 1 - \left(\frac{\mathrm{SSE}}{\mathrm{TSS}} \cdot \frac{n-1}{n-k-1}\right)$$

• This adjusted R^2 is usually a bit smaller than R^2 , and the difference decreases as n gets large.

Multiple Linear Regression 000000	Model Assumptions for MLR 0000	Testing Model Fit 00000	Model Selection

• Consider a (hypothetical) statistics class of 30 students; suppose these students have two midterm exams and a final.

Multiple Linear Regression 000000	Model Assumptions for MLR	Testing Model Fit 00000	Model Selection

- Consider a (hypothetical) statistics class of 30 students; suppose these students have two midterm exams and a final.
 - Can we use student scores on the midterms to predict their score on the final?

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit 00000	Model Selection

- Consider a (hypothetical) statistics class of 30 students; suppose these students have two midterm exams and a final.
 - Can we use student scores on the midterms to predict their score on the final?



Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit 00000	Model Selection

- Consider a (hypothetical) statistics class of 30 students; suppose these students have two midterm exams and a final.
 - Can we use student scores on the midterms to predict their score on the final?



• Individually, both exams seem to have relatively strong linear relationship with the final.

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	Model Selection
			00000000

Checking Conditions

• Are conditions for inference met?

Multiple	Linear	Regression
00000		

Model Assumptions for MLR 0000 Testing Model Fit 00000 Model Selection

Checking Conditions

• Are conditions for inference met?



Multiple	Linear	Regression

Model Assumptions for MLR 0000 Testing Model Fit 00000 Model Selection

Checking Conditions

• Are conditions for inference met?



- Linearity
- Ø Independence
- 8 Normality
- 4 Equal Variability

	Model Assumptions for MLR	Testing Model Fit	Model Selection
000000	0000	00000	00000000

• Let's build the model:

	Model Assumptions for MLR	Testing Model Fit	Model Selection
000000	0000	00000	00000000

• Let's build the model: scores_mod <- lm(final ~ exam1 + exam2, data = stat_scores) get_regression_table(scores_mod)

##		term	estimate	std_error	statistic	p_value	lower_ci	upper_ci
##	1	intercept	4.830	9.76	0.49	0.625	-15.20	24.86
##	2	exam1	0.836	0.27	3.11	0.004	0.28	1.39
##	3	exam2	0.099	0.25	0.40	0.694	-0.41	0.61

Model Assumptions for MLR	Testing Model Fit	Model Selection
		00000000

• Let's build the model: scores_mod <- lm(final ~ exam1 + exam2, data = stat_scores) get_regression_table(scores_mod)

##		term	estimate	std_error	statistic	p_value	lower_ci	upper_ci
##	1	intercept	4.830	9.76	0.49	0.625	-15.20	24.86
##	2	exam1	0.836	0.27	3.11	0.004	0.28	1.39
##	3	exam2	0.099	0.25	0.40	0.694	-0.41	0.61

• Based on individual slope hypothesis tests...

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	Model Selection
			00000000

• Let's build the model: scores_mod <- lm(final ~ exam1 + exam2, data = stat_scores) get_regression_table(scores_mod)

##		term	estimate	std_error	statistic	p_value	lower_ci	upper_ci
##	1	intercept	4.830	9.76	0.49	0.625	-15.20	24.86
##	2	exam1	0.836	0.27	3.11	0.004	0.28	1.39
##	3	exam2	0.099	0.25	0.40	0.694	-0.41	0.61

• Based on individual slope hypothesis tests...

• Exam1 is a significant predictor of final score, while Exam2 is **not** a significant predictor.

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	Model Selection
			00000000

• Let's build the model: scores_mod <- lm(final ~ exam1 + exam2, data = stat_scores) get_regression_table(scores_mod)

##		term	estimate	std_error	statistic	p_value	lower_ci	upper_ci
##	1	intercept	4.830	9.76	0.49	0.625	-15.20	24.86
##	2	exam1	0.836	0.27	3.11	0.004	0.28	1.39
##	3	exam2	0.099	0.25	0.40	0.694	-0.41	0.61

Based on individual slope hypothesis tests...

- Exam1 is a significant predictor of final score, while Exam2 is not a significant predictor.
- Hang on... let's build a simple linear model for final and exam2

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	Model Selection
			00000000

• Let's build the model: scores_mod <- lm(final ~ exam1 + exam2, data = stat_scores) get_regression_table(scores_mod)

##		term	estimate	std_error	statistic	p_value	lower_ci	upper_ci
##	1	intercept	4.830	9.76	0.49	0.625	-15.20	24.86
##	2	exam1	0.836	0.27	3.11	0.004	0.28	1.39
##	3	exam2	0.099	0.25	0.40	0.694	-0.41	0.61

Based on individual slope hypothesis tests...

• Exam1 is a significant predictor of final score, while Exam2 is **not** a significant predictor.

• Hang on... let's build a simple linear model for final and exam2

```
scores_mod2 <- lm(final ~ exam2, data = stat_scores)
get_regression_table(scores_mod2)</pre>
```

##		term	estimate	std_error	statistic	p_value	lower_ci	upper_ci
##	1	intercept	16.86	10.26	1.6	0.11	-4.2	37.9
##	2	exam2	0.78	0.14	5.7	0.00	0.5	1.1

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	Model Selection
			00000000

• Let's build the model: scores_mod <- lm(final ~ exam1 + exam2, data = stat_scores) get_regression_table(scores_mod)

##		term	estimate	std_error	statistic	p_value	lower_ci	upper_ci
##	1	intercept	4.830	9.76	0.49	0.625	-15.20	24.86
##	2	exam1	0.836	0.27	3.11	0.004	0.28	1.39
##	3	exam2	0.099	0.25	0.40	0.694	-0.41	0.61

Based on individual slope hypothesis tests...

• Exam1 is a significant predictor of final score, while Exam2 is **not** a significant predictor.

 Hang on... let's build a simple linear model for final and exam2 scores mod2 <- lm(final - exam2, data = stat scores)

```
get_regression_table(scores_mod2)
```

##		term	estimate	std_error	statistic	p_value	lower_ci	upper_ci
##	1	intercept	16.86	10.26	1.6	0.11	-4.2	37.9
##	2	exam2	0.78	0.14	5.7	0.00	0.5	1.1

In the simple mode, exam 2 is a significant predictor of final.

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	Model Selection
000000	0000	00000	00000●00

• What is going on here? (Note that exam2 had a correlation of 0.74 with the final)

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	Model Selection
000000	0000	00000	

- What is going on here? (Note that exam2 had a correlation of 0.74 with the final)
 - Recall the form of the individual slope hypotheses:

 $H_0: \beta_{exam2} = 0$, given that exam1 is in the model

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	Model Selection
			00000000

- What is going on here? (Note that exam2 had a correlation of 0.74 with the final)
 - Recall the form of the individual slope hypotheses:

 $H_0: \beta_{exam2} = 0$, given that exam1 is in the model

• On its own, Exam 2 is helpful for predicting the final. But IF Exam 1 is already in the model, Exam 2 is redundant.

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	Model Selection
			00000000

- What is going on here? (Note that exam2 had a correlation of 0.74 with the final)
 - Recall the form of the individual slope hypotheses:

 $H_0: \beta_{exam2} = 0$, given that exam1 is in the model

- On its own, Exam 2 is helpful for predicting the final. But IF Exam 1 is already in the model, Exam 2 is redundant.
- Note that Exam 1 and Exam 2 are highly correlated



Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	Model Selection
			00000000

- What is going on here? (Note that exam2 had a correlation of 0.74 with the final)
 - Recall the form of the individual slope hypotheses:

 $H_0: \beta_{exam2} = 0$, given that exam1 is in the model

- On its own, Exam 2 is helpful for predicting the final. But IF Exam 1 is already in the model, Exam 2 is redundant.
- Note that Exam 1 and Exam 2 are highly correlated



• Once exam 1 is known, exam 2 doesn't contribute much additional information.

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	Model Selection
000000	0000	00000	

Model Selection

• We have 3 models for predicting final exam score:

 ${\rm final} \sim {\rm exam1} \quad {\rm final} \sim {\rm exam2} \quad {\rm final} \sim {\rm exam1} + {\rm exam2}$

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	Model Selection	
000000	0000	00000	000000●0	

Model Selection

• We have 3 models for predicting final exam score:

```
final \sim exam1 final \sim exam2 final \sim exam1 + exam2
```

• To decide which is best, let's perform F tests and calculate R^2 and R^2_{adj} :

```
get_regression_summaries(scores_mod1)
```

A tibble: 1 x 9
r_squared adj_r_squared mse rmse sigma statistic p_value df nobs
<dbl> <dbl > <

A tibble: 1 x 9
r_squared adj_r_squared mse rmse sigma statistic p_value df nobs
<dbl> <dbl > <

##	#	A tibble:	1 x 9							
##		r_squared	adj_r_squared	mse	rmse	sigma	statistic	p_value	df	nobs
##		<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>
##	1	0.662	0.637	35.3	5.94	6.26	26.5	0	2	30

	Model Assumptions for MLR	Testing Model Fit	Model Selection
000000	0000	00000	00000000

• All 3 models had statistically significant F-statistics

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	Model Selection
			0000000

- All 3 models had statistically significant F-statistics
- Model 1 had $R^2 = 0.66$, Model 2 had $R^2 = 0.541$, and the full Model had $R^2 = 0.662$

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	Model Selection
			0000000

- All 3 models had statistically significant F-statistics
- Model 1 had $R^2 = 0.66$, Model 2 had $R^2 = 0.541$, and the full Model had $R^2 = 0.662$
- But Model 1 had $R^2_{\rm adj}=$ 0.648, Model 2 had $R^2_{\rm adj}=$ 0.525, and the full Model had $R^2_{\rm adj}=$ 0.637

Multiple Linear Regression	Model Assumptions for MLR	Testing Model Fit	Model Selection
			0000000

- All 3 models had statistically significant F-statistics
- Model 1 had $R^2 = 0.66$, Model 2 had $R^2 = 0.541$, and the full Model had $R^2 = 0.662$
- But Model 1 had $R^2_{\rm adj}=$ 0.648, Model 2 had $R^2_{\rm adj}=$ 0.525, and the full Model had $R^2_{\rm adj}=$ 0.637
- Since Model 1 had the highest adjusted R^2 and was the simplest model considered, Model 1 is likely the most accurate of the 3 models.